

Bayesian learning Methods

- Bayesian learning Methods are relevant to our study of machine learning for two different reasons:
- ① Bayesian learning algorithms that calculate explicit probabilities for hypotheses, such as Naive Bayes Classifier, are among the most practical approaches to certain types of learning problems.
 - ② The second reason is that Bayesian methods provide a useful perspective for understanding many learning algorithms that do not explicitly manipulate probabilities.

Two rules for Bayesian methods:

- ③ Provide practical learning algorithms

- Naive Bayes learning → each observed training example can incrementally decrease or increase the estimated probability that a hypothesis is correct. This provides a more flexible approach to learning than algorithms that completely eliminate a hypothesis if it is found to be inconsistent with any single example.
 - Combine prior knowledge (prior probabilities) with observed data;
→ prior knowledge can be combined with observed data to determine the final probability of a hypothesis;
 - Make probabilistic predictions
→ Bayesian methods can accommodate hypotheses that make probabilistic predictions. New instances can be classified by combining the predictions of multiple hypotheses, weighted by their probabilities;
- ② Provide useful conceptual framework
- Provide a standard for evaluating other learning algorithms
→ even in cases where Bayesian methods prove computationally intractable, they can provide a standard of optimal decision making against which other practical methods can be measured

One practical difficulty in applying Bayesian methods is that they typically require initial knowledge of many probabilities. When these probabilities are not known in advance they are often estimated based on background knowledge, previously available data and assumptions about the form of the underlying distributions. A second, practical difficulty is the significant computational cost required to determine the Bayes optimal hypothesis in the general case.

Basic formulas for probabilities

Product rule : probability of conjunction of A and B

$$P(A \wedge B) = P(A|B)P(B) = P(B|A)P(A)$$

Sum rule : probability of disjunction of A and B

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

Theorem of total probability : if events A_1, \dots, A_n are mutually exclusive with $\sum_{i=1}^n P(A_i) = 1$, then

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

Bayes theorem : provides a way to calculate posterior probability $P(h|D)$ from the prior probability $P(h)$ together with $P(D)$ and $P(D|h)$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

In machine learning we are often interested in determining the best hypothesis from some space H , given the observed training data D . One way to specify what we mean by the best hypothesis is to say that we demand the most probable one, given the data D plus any initial knowledge about the prior probabilities of the various hypotheses in H . Bayes theorem provides a direct method for computing such probabilities.

We want to learn a function with a dataset with the goal of classifying a new instance. Given a target function $f: X \rightarrow V$, a dataset D and a new instance $x' \notin D$, the best prediction $f(x') = v^*$ is:

$$v^* = \operatorname{argmax}_{v \in V} P(v | x', D)$$

In other words, given the dataset D and a new sample $x' \notin D$, we want to compute the probability distribution over V : for each possible outcome of $f: X \rightarrow V$, we compute the likelihood that it is the right one and take the sample with the higher probability.

$$P(V | x', D) \leftarrow \text{probability distribution on the image of } f(V)$$

Given a dataset D and hypothesis space H , compute a probability distribution over H given D .

$$P(H | D)$$

We could use the Bayes rule if H is finite and small

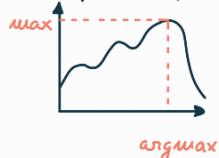
$$P(h | D) = \frac{P(D|h) P(h)}{P(D)}$$

↑
for each h

$P(h)$ = prior probability of hypothesis h
 $P(D)$ = prior probability of training data D
 $P(h | D)$ = probability of h given D - unknown
 $P(D|h)$ = probability of D given h

Maximum A Posteriori Hypothesis h_{MAP}

In many learning scenarios, the learner consider some set of candidate hypotheses H and is interested in finding the most probable hypothesis $h \in H$ given the observed data D . Any such maximally probable hypothesis is called maximum a posteriori (MAP) hypothesis. We can determine the MAP hypotheses by using Bayes theorem to compute the posterior probability for each candidate.



$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h | D) = \operatorname{argmax}_{h \in H} \frac{P(D|h) P(h)}{P(D)} = \operatorname{argmax}_{h \in H} P(D|h) P(h)$$

argmax is invariant with respect to a scaling factor ($\frac{1}{P(D)}$) which is constant with respect to h

Maximum Likelihood Hypotheses

$P(D|h)$ is often called the likelihood of the data D given h and any hypothesis that maximizes $P(D|h)$ is called a maximum likelihood hypothesis (ML hypothesis). We don't have info about prior probabilities of the hypotheses so we assume they all have the same prior probability; the hypothesis space is uniformly distributed. If we assume $P(h_i) = P(h_j) \forall i, j$ then we can further simplify and choose the maximum likelihood (ML) hypothesis

$$h_{ML} = \operatorname{argmax}_{h \in H} P(D|h)$$

We can design a straightforward concept learning algorithm to output the maximum a posteriori hypothesis, based on the Bayes theorem:

■ for each hypothesis $h \in H$, compute the posterior probability with Bayes $P(h | D) = \frac{P(D|h) P(h)}{P(D)} \quad \forall h$

■ output the hypothesis h_{MAP} with the highest posterior probability $h_{MAP} = \operatorname{argmax}_{h \in H} P(h | D)$

We can implement this algorithm when it is feasible to enumerate every h in H .

This algorithm requires significant amounts of computation because it applies the Bayes theorem to each hypothesis h in H to compute $P(h | D)$. While this may prove impractical for large hypothesis spaces, the algorithm is still of interest because it provides a standard against which we may judge the performance of other concept learning algorithms.

Most probable classification \neq Most probable hypothesis

h_{MAP} = most probable hypothesis given data D

Given a new instance x , what is its most probable classification?

$h_{MAP}(x)$ may not be the most probable classification

Consider three possible hypotheses h_1, h_2, h_3

$$P(h_1|D) = 0.6, P(h_2|D) = 0.3, P(h_3|D) = 0.3$$

Given a new instance x

$$h_1(x) = \oplus, h_2(x) = \ominus, h_3(x) = \ominus$$

x has probability 0.6 (40%) of being classified as a \oplus and

The most probable classification is \ominus (60%)

The maximum a posteriori hypothesis is $h_{MAP} = \arg\max_{h_i \in H} P(h_i|D)$

$$\begin{aligned} h_{MAP} &= \arg\max_{h_i \in H} P(h_i|D) = \\ &= \arg\max_{h_i \in H} P(D|h_i)P(h_i) \\ P(h_i) \text{ is constant} &\rightarrow \arg\max_{h_i \in H} P(D|h_i) \\ 0.3 + 0.3 = 0.6 &\text{ of being classified as } \ominus \end{aligned}$$

Bayes Optimal Classifier

Consider a target function $f: X \rightarrow V, V = \{v_1, \dots, v_k\}$, data set D

and a new instance $x \notin D$

$$P(v_j|x, D) = \sum_{h_i \in H} P(v_j|x, h_i) P(h_i|D)$$

x is not in D (h_i is conditionally independent from D)

$$= \sum_{h_i \in H} P(v_j|x, h_i) P(h_i|D)$$

once we know h_i we don't need the dataset
since it is derived from it

each $P(v_j|x, h_i)$ is scaled
down by $P(h_i|D)$
 $\sum P(h_i|D) = 1$

$$BOC \rightarrow V_{OB} = \arg\max_{v_j \in V} \sum_{h_i \in H} P(v_j|x, h_i) P(h_i|D)$$

Consider $H = \{h_1, h_2, h_3\}$

$$\begin{array}{lll} h_{MAP} & P(h_1|D) = 0.4 & P(\ominus|x, h_1) = 0 \\ & P(h_2|D) = 0.3 & P(\oplus|x, h_1) = 1 \\ & P(h_3|D) = 0.3 & P(\ominus|x, h_2) = 1 \\ & & P(\oplus|x, h_2) = 0 \\ & & P(\ominus|x, h_3) = 1 \\ & & P(\oplus|x, h_3) = 0 \end{array}$$

Suppose a new instance x is encountered and classified as \oplus by h_1 and as \ominus by h_2, h_3

$$\sum_{h_i \in H} P(\oplus|x, h_i) P(h_i|D) = 0.4$$

the most probable classification (\ominus) is different

$$\sum_{h_i \in H} P(\ominus|x, h_i) P(h_i|D) = 0.6 \rightarrow$$

from the classification generated by $h_{MAP}(h_1)$

The most probable classification is obtained by combining the predictions of all hypotheses, weighting them by the posterior probability of the generating hypothesis

$$V_{OB} = \arg\max_{v_j \in V} \sum_{h_i \in H} P(v_j|x, h_i) P(h_i|D) = \ominus$$

This method is good but it is not applicable because we can't iterate over all hypotheses.
We need another method \rightarrow BOC is optimal but not applicable when H is huge or infinite

If $P(D|h_i)$ is a continuous function (continuous space of hypothesis leads to continuous probability distribution) we can find the argmax of $\log(P)$

$$h_{ML} = \arg\max_{h_i \in H} P(D|h_i) = \arg\max_{h_i \in H} (\log(P(D|h_i)))$$

The log turns the product (conditional independence $P(D|h) = \prod_i P(d_i|h_i)$) into sum

We have 5 kinds of bags containing 2 kinds of candies (l, c) in different proportions

- 20% of the bags are b_1 : 100% c , 0% l
- 20% of the bags are b_2 : 75% c , 25% l
- 20% of the bags are b_3 : 50% c , 50% l
- 20% of the bags are b_4 : 25% c , 75% l
- 10% of the bags are b_5 : 0% c , 100% l

We choose a random bag (without knowing which type it is) and extract some candies from it.

What kind of bag is it?

What is the probability of extracting a candy of a specific flavor next?

$$\begin{array}{ll} \text{Prior probability distribution: } & P(H) = \langle 0.1, 0.2, 0.4, 0.2, 0.1 \rangle \\ \text{likelihood for } l \text{ candy} & P(l|H) = \langle 0, 0.25, 0.5, 0.75, 1 \rangle \\ \text{likelihood for } c \text{ candy} & P(c|H) = \langle 1, 0.75, 0.5, 0.25, 0 \rangle \end{array}$$

What kind of bag is it?

$$P(l_i | D) = \propto P(D | l_i) P(l_i) \quad \leftarrow \text{Bayes rule}$$

We need a dataset. Suppose we extract 1 l candy

$$\left. \begin{array}{l} P(l_1 | d_1) = \propto P(d_1 | l_1) P(l_1) \\ = \propto P(l | l_1) P(l_1) \\ = \propto \cdot 0 \cdot 0.1 \\ P(l_2 | d_1) = \propto \cdot 0.25 \cdot 0.2 \\ P(l_3 | d_1) = \propto \cdot 0.5 \cdot 0.4 \\ P(l_4 | d_1) = \propto \cdot 0.75 \cdot 0.2 \\ P(l_5 | d_1) = \propto \cdot 1 \cdot 0.1 \end{array} \right\} \begin{array}{l} \text{only one extraction} \\ \uparrow \\ P(H | D = \{d_1\}) = \propto \cdot \langle 0, 0.25, 0.5, 0.75, 1 \rangle \cdot \langle 0.1, 0.2, 0.4, 0.2, 0.1 \rangle \\ = \propto \langle 0, 0.05, 0.2, 0.15, 0.1 \rangle \end{array}$$

Suppose we extract another l candy

$D_2 = \{d_1, d_2\}$ subsequent extractions are independent from each other
Bayes + chain rule

$$\begin{aligned} P(l_i | \{d_1, d_2\}) &= \propto P(\{d_1, d_2\} | l_i) P(l_i) \\ &= \propto P(\{d_2\} | l_i) \underbrace{P(\{d_1\} | l_i) P(l_i)}_{\text{previous step}} \\ &= \propto \cdot \langle 0, 0.25, 0.05, 0.75, 1 \rangle \cdot \propto \cdot \langle 0, 0.05, 0.2, 0.15, 0.1 \rangle = \dots \end{aligned}$$

What is the probability of having another candy l after $D_3 = \{d_1, d_2, d_3\} = ?$

$$P(l | D_3) = \sum_{l_i \in H} \underbrace{P(l | l_i)}_{\text{vote of each hypothesis for } l} \underbrace{P(l_i | D_3)}_{\text{weights of each hypothesis (how likely this hypothesis is)}} = \text{optimal prediction}$$

General approach
 Given dataset $D = \{d_i\}$ with $d_i \in \{0, 1\}$
 assuming a probability distribution $P(d_i, \theta)$ → parametric probability distribution parametrized by θ

$$\theta_{\text{re}} = \underset{\theta}{\operatorname{argmax}} \log P(d_i | \theta)$$

we can apply the same model to multiple distributions (e.g. Bernoulli)

Bernoulli distribution is the discrete probability distribution of a binary random variable $X \in \{0, 1\}$
 $P(X=1) = \theta$; $P(X=0) = 1-\theta$

(e.g. observing head after flipping a coin, extracting candy l from a set containing 2 candies, ...)

$$P(X=k, \theta) = \theta^k (1-\theta)^{1-k}$$

Multivariate Bernoulli distribution is the joint probability distribution of a set of binary random variables X_1, \dots, X_m , each of them following the Bernoulli distribution.

this under the assumption that random variables X_i are mutually independent → Multivariate Bernoulli is the product of m Bernoulli distributions

e.g. observing head after flipping a coin AND extracting candy l

$$X_1 = \text{head} \quad X_2 = l$$

$$P(X_1=k, X_2=l, \theta_1, \theta_2) = \theta_1^{k_1} (1-\theta_1)^{1-k_1} \cdot \theta_2^{k_2} (1-\theta_2)^{1-k_2}$$

θ_1 = number of head observed / number of coin flip $\sim \frac{1}{2}$

θ_2 = number of l candies / number of candies in the bag $\in [0, 1]$

Binomial distribution is the probability distribution of k outcomes from n Bernoulli trials

$$P(X=k, n, \theta) = \binom{n}{k} \theta^k (1-\theta)^{n-k} \text{ Bernoulli}$$

\uparrow k outcomes over n trials

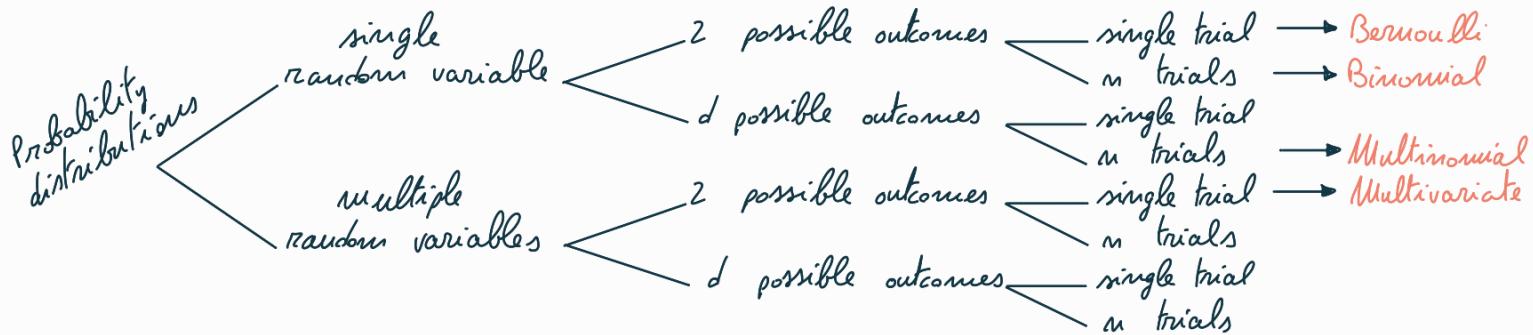
e.g. flipping a coin n times and observing k heads

Multinomial distribution is the generalization of binomial distribution for discrete valued random variables with d possible outcomes.

X_d has d possible outcomes (k_1, \dots, k_d). What is the probability of observing a particular outcome k times after n trials?

$$P(X_1=k_1, \dots, X_d=k_d, n, \theta_1, \dots, \theta_d) = \frac{n!}{k_1! \dots k_d!} \theta_1^{k_1} \dots \theta_d^{k_d}$$

e.g. rolling a d -sided die n times and observing k times a particular value
 extracting k candies after n trials from a bag containing d different flavors



Bernoulli	: single random variable	, binary	, single trial
Binomial	: single random variable	, binary	, n trials
Multinomial	: single random variable	, multinomial	, n trials
Multivariate	: multiple random variables	, binary	, single trial

Probabilistic classification : $\operatorname{argmax}_{v_j \in V} P(v_j | x, D)$

Bayes Optimal Classifier provides best result, not practical when hypothesis space is large

Continuous model

Maximum likelihood estimation efficiently solved when analytical solutions are available

Naive Bayes Classifier is an approximation of the Bayes Optimal Classifier. It uses conditional independence to approximate the solution. It is used when each instance is described by a set of attributes : $f : A_1 \times A_2 \times \dots \times A_m \rightarrow V = \{v_1, v_2, \dots, v_k\}$

$$P(x, y | z) = P(x | y, z)P(y | z) = P(x | z)P(y | z) \quad + \quad \operatorname{argmax}_{v_j \in V} P(v_j | x, D) = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2, \dots, a_m, D)$$

↑ the same without y

x is conditionally independent of y given z

BC with x described as a set of features

$$v_{\text{MAP}} = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2, \dots, a_m, D) = \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2, \dots, a_m | v_j, D) P(v_j | D)}{P(a_1, a_2, \dots, a_m | D)} =$$

normalization factor

$$= \operatorname{argmax}_{v_j \in V} P(a_1, a_2, \dots, a_m | v_j, D) P(v_j | D) = \operatorname{argmax}_{v_j \in V} P(v_j | D) \prod_i P(a_i | v_j, D)$$

conditional independence assumption

The Naive Bayes assumes that all attributes are mutually conditionally independent from each other given the dataset and the classification value; this is not real in general but we are approximating the solution; it is therefore no more guaranteed to get the optimal value. If we don't make this assumption the joint probability $P(a_1, a_2, \dots, a_m)$ will be exponential in the number of attributes and will make the system impractical. Conditional independence can bring exponential systems to linear.

$$v_{\text{NB}} = \operatorname{argmax}_{v_j \in V} P(v_j | D) \prod_i P(a_i | v_j, D)$$

implicit dependency on the dataset

This is a practical method

Naive Bayes Algorithm

$$f : A_1 \times A_2 \times \dots \times A_m \rightarrow V = \{v_1, v_2, \dots, v_k\}$$

dataset D

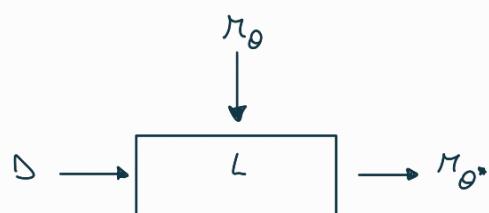
new instance $x = \langle a_1, a_2, \dots, a_m \rangle$

}

inputs

Naive-Bayes-learner (A, V, D) :

for each target value $v_j \in V$
 $\hat{P}(v_j | D) \leftarrow \text{estimate } P(v_j | D)$
 for each attribute $A_i \in A$
 for each attribute value $a_i \in A_k$
 $\hat{P}(a_i | v_j, D) \leftarrow \text{estimate } P(a_i | v_j, D)$

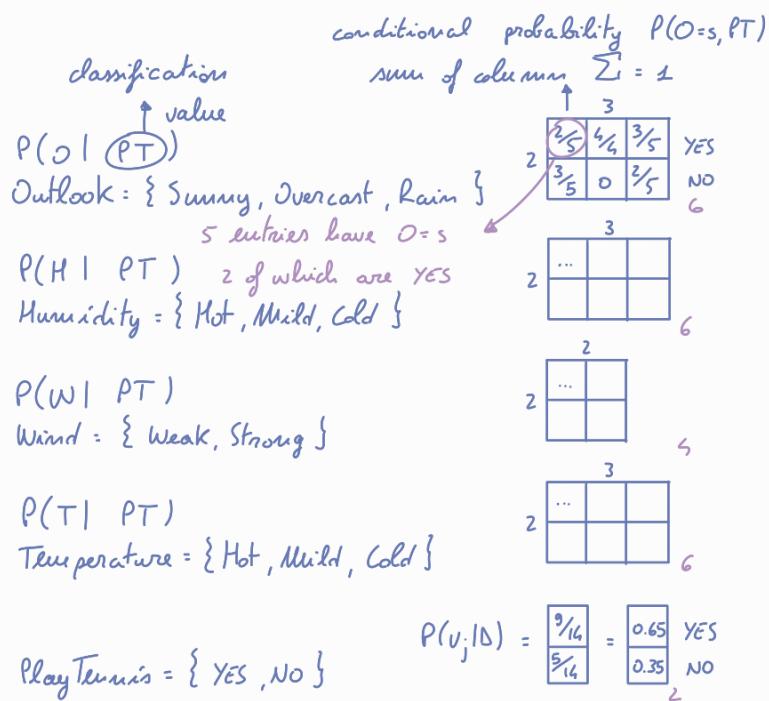


now we have $\hat{P}(v_j | D)$ and $\hat{P}(a_i | v_j, D)$ for each attribute value for each attribute

Classify-new-instance (x) : $v_{\text{NB}} = \operatorname{argmax}_{v_j \in V} \hat{P}(v_j | D) \prod_{a_i \in x} \hat{P}(a_i | v_j, D)$

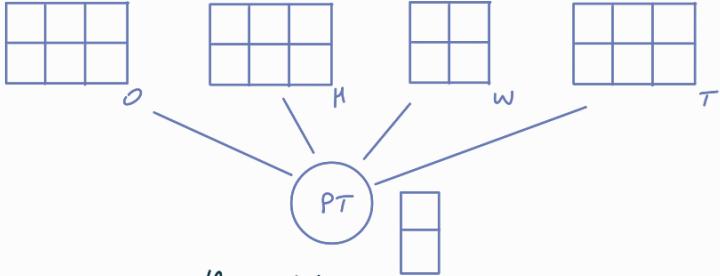
e.g. playtennis example

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No



The model will be complete when we'll have these $6+6+6+6+2 = 26$ numbers
 We didn't account for the fact that the sum of each column is 1. The number of trainable parameters (size of the model) is exactly half since we only have the freedom to choose 1 row (the other one is complementary). This is a special case, not always the size is exactly half.
 Size of the model = $26/2 = 13$

We can consider the problem as a Bayesian network



What if a $P(a_i|v_j, \Delta)$ is 0?

$$P(A_2 = a_2 | v_j, \Delta) = 0 \longrightarrow \prod_i P(\cdot) = 0$$

This happens if an attribute is never observed in the dataset

e.g. never observed in the dataset that the output is NO when Outlook is Overcast

Typical solution is Bayesian estimate with prior estimates

$$\hat{P}(v_j|\Delta) = \frac{|\{ \dots, v_j \}|}{|\Delta|} = \frac{\text{times class } v_j \text{ is observed in } \Delta}{\text{number of elements in } \Delta}$$

$$\hat{P}(a_i|v_j, \Delta) = \frac{|\{ \dots, a_i, \dots, v_j \}|}{|\{ \dots, v_j \}|} = \frac{\text{times attribute } A_i \text{ has value } a_i \text{ when class } v_j \text{ is predicted}}{\text{times class } v_j \text{ is observed in } \Delta}$$

$$\frac{|\{ \dots, a_i, \dots, v_j \}| + m p}{|\{ \dots, v_j \}| + m}$$

p is a prior estimate for $P(a_i|v_j, \Delta)$

$$m$$
 is a weight given to prior

new instance $x = \langle \text{Outlook} = s, \text{Temperature} = c, \text{Humidity} = h, \text{Wind} = w \rangle$

$$V_{NB} = \underset{v_j \in \{\text{YES}, \text{NO}\}}{\operatorname{argmax}} P(v_j) \prod_i P(a_i|v_j)$$

$$= \underset{v_j \in \{\text{YES}, \text{NO}\}}{\operatorname{argmax}} P(v_j) P(\text{Outlook} = s|v_j) P(\text{Temperature} = c|v_j) P(\text{Humidity} = h|v_j) P(\text{Wind} = w|v_j)$$

Conditional independence assumption is often violated

N.B.: conditioning on Δ is omitted